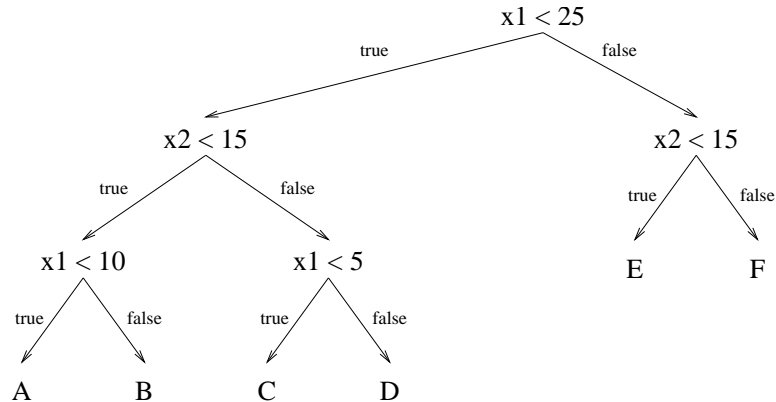


AI534 — Written Homework Assignment 3 —

1. (6 pts) Consider the following decision tree:



- (a) (2 pts) Draw the decision boundaries defined by this tree. Each leaf of the tree is labeled with a letter. Write this letter in the corresponding region of input space.
- (b) (2 pts) Give another decision tree that is syntactically different but defines the same decision boundaries. This demonstrates that the space of decision trees is syntactically redundant.
- (c) (2pts) How does this redundancy influence learning (does it make it easier or harder to find an accurate tree)?
2. (6 pts) In the basic decision tree algorithm (assuming we always create binary splits), we choose the feature/value pair with the maximum information gain as the test to use at each internal node of the decision tree. Suppose we modified the algorithm to choose at random from among those feature/value combinations that had non-zero mutual information, and we kept all other parts of the algorithm unchanged.
- (a) (2 pts) What is the maximum number of leaf nodes that such a decision tree could contain if it were trained on m training examples?
- (b) (2 pts) What is the maximum number of leaf nodes that a decision tree could contain if it were trained on m training examples using the original maximum mutual information version of the algorithm? Is it bigger, smaller, or the same as your answer to (b)?
- (c) (2 pts) How do you think this change (using random splits vs. maximum information mutual information splits) would affect the testing accuracy of the decision trees produced on average? Why?
3. (8 pts) Consider the following training set:

A	B	C	Y
0	1	1	0
1	1	1	0
0	0	0	0
1	1	0	1
0	1	0	1
1	0	1	1

Learn a decision tree from the training set shown above using the information gain criterion. Show your steps, including the calculation of information gain (you can skip $H(y)$ and just compute $H(y|\mathbf{x})$) of different candidate tests. You can randomly break ties (or better, choose the one that give you smaller tree if you do a bit look ahead for this problem).

4. (5 pts) Please show that in iteration l of Adaboost, the weighted error of h_l on the updated weights D_{l+1} is exactly 50%. In other words, $\sum_{i=1}^N D_{l+1}(i) I(h_l(X_i) \neq y_i) = 50\%$, where $I(\cdot)$ is the indicator function that takes value 1 if the argument is true. (Hint: start with the condition that, post update, the total weight of correct examples equals the total weight of in-correct examples, i.e., 50% each.)
5. **HAC.** (4pts) Create by hand the clustering dendrogram for the following samples of ten points in one dimension.

$$Sample = (-2.2, -2.0, -0.3, 0.1, 0.2, 0.4, 1.6, 1.7, 1.9, 2.0)$$

- a. (2pts) Using single link.
 - b. (2pts) Using complete link
6. (6 pts) Deriving Kmeans for L_1 norm. Consider replacing the distance function used for Kmeans with L_1 norm with the following objective:

$$\min_{\mu_1, \dots, \mu_K, C_1, \dots, C_K} \sum_{i=1}^K \sum_{x \in C_i} |x - \mu_i|$$

- (3 pts) Show that given fixed cluster assignments C_1, \dots, C_K , the prototype μ_i that optimizes the above objective can be obtained by taking element wise median of all the points in cluster i .
- (3 pts) Modify the kmeans algorithm for this L_1 based objective.